

GNNs for Recommendation Systems

Prof. O-Joun Lee

Dept. of Artificial Intelligence,
The Catholic University of Korea
ojlee@catholic.ac.kr

Contents



- Tasks and Evaluation metrics
- Embedding models
- GNNs for Recommendation systems
 - Traditional Collaborative Filtering
 - Neural Graph Collab. Filtering (NGCF)
 - LightGCN
 - Recommender Systems With Knowledge Graph

- **Information Explosion in the era of Internet**

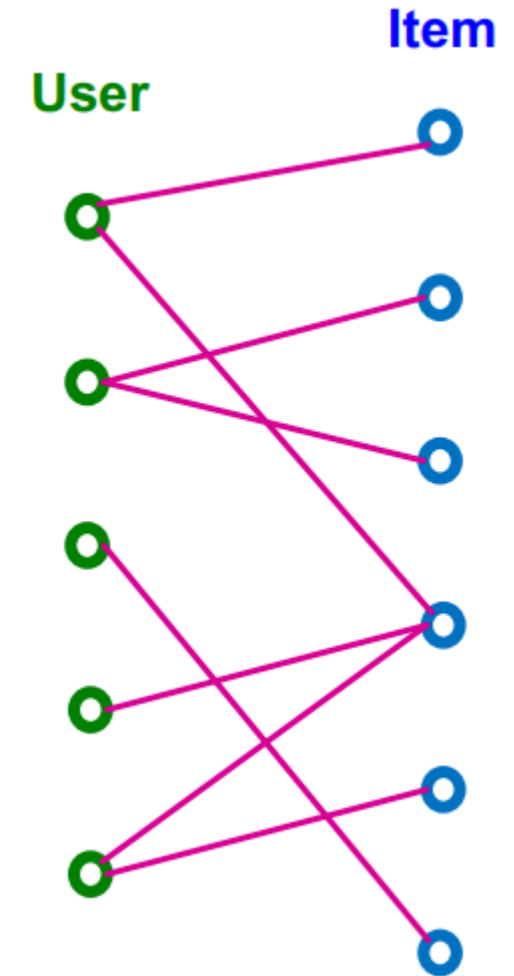
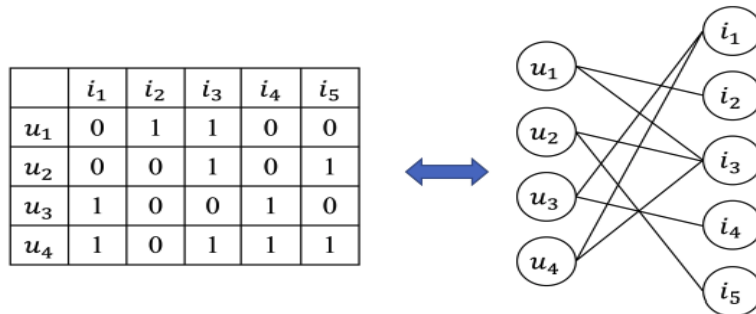
- 10K+ movies in Netflix
- 12M products in Amazon (350m on Marketplace)
- 70M+ music tracks in Spotify
- 10B+ videos on YouTube
- 200B+ pins (images) in Pinterest

- **Personalized recommendation**

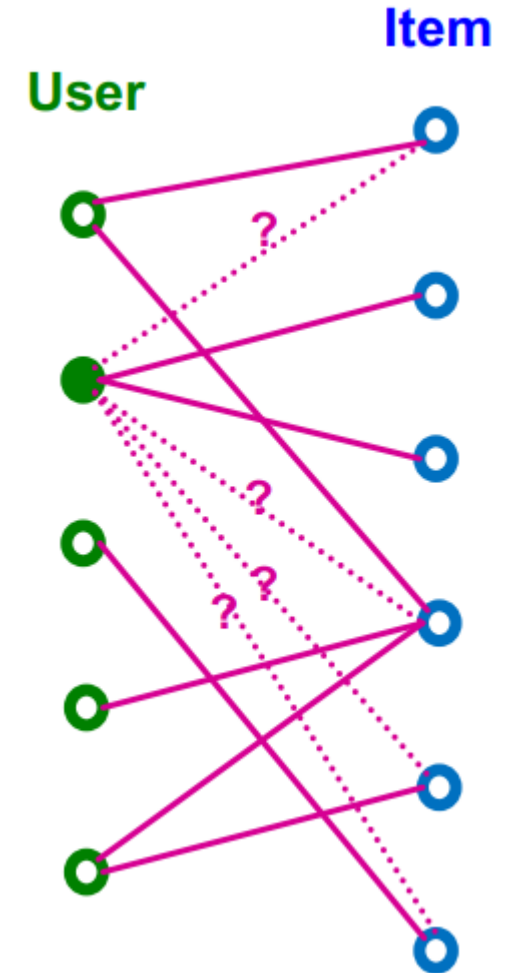
- i.e., suggesting a small number of interesting items for each user) is critical for users to effectively explore the content of their interest.

- **What are the inputs of recommender system?**
 - Main info:
 - user (attributes)
 - item (attributes)
 - user feedback on items (user-item interactions)
 - Side info: (to solve data sparsity and cold start issue)
 - social relationship between users
 - knowledge graph
- **Tasks:**
 - users' preferences
 - item properties

- Recommender system can be naturally modeled as a **bipartite graph**
 - A graph with two node types: users and items.
 - **Edges** connect **users** and **items**
 - Indicates user-item interaction (e.g., click, purchase, review etc.)
 - Often associated with timestamp (timing of the interaction).

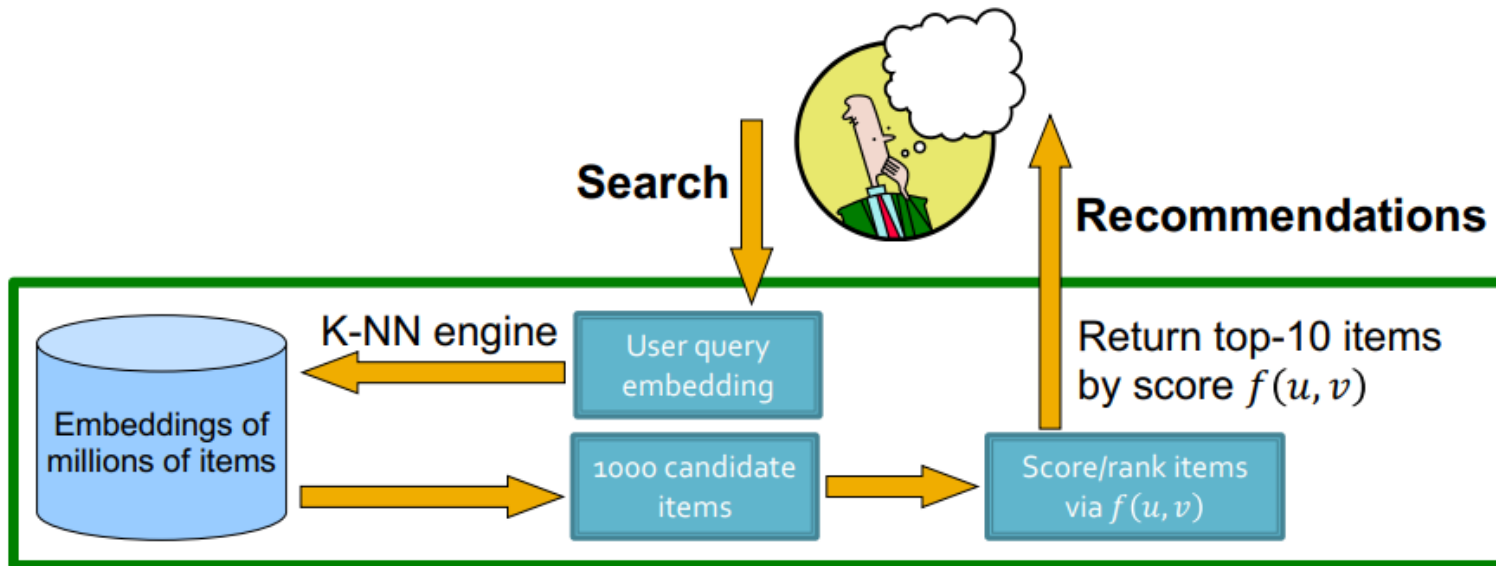


- **Given:**
 - Past user-item interactions
- **Task:**
 - Predict new items each user will interact in the future.
 - Can be cast as link prediction problem.
 - Predict new user-item interaction edges given the past edges.
 - For $u \in U, v \in V$, we need to get a real-valued score $f(u, v)$



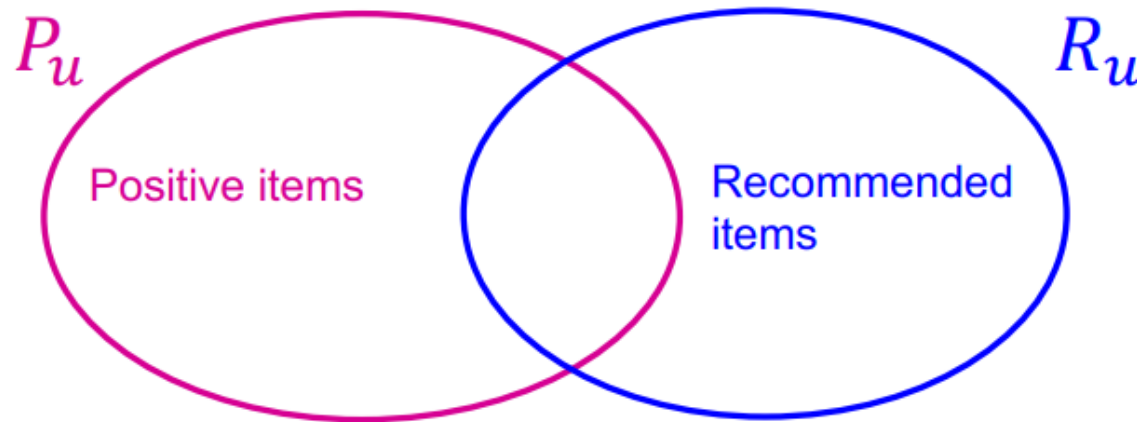
- Problem: Cannot evaluate $f(u, v)$ for every user u – item v pair.
- Solution: 2-stage process:
 - Candidate generation (cheap, fast)
 - Ranking (slow, accurate)

Example $f(u, v)$:
 $f(u, v) = z_u \cdot z_v$

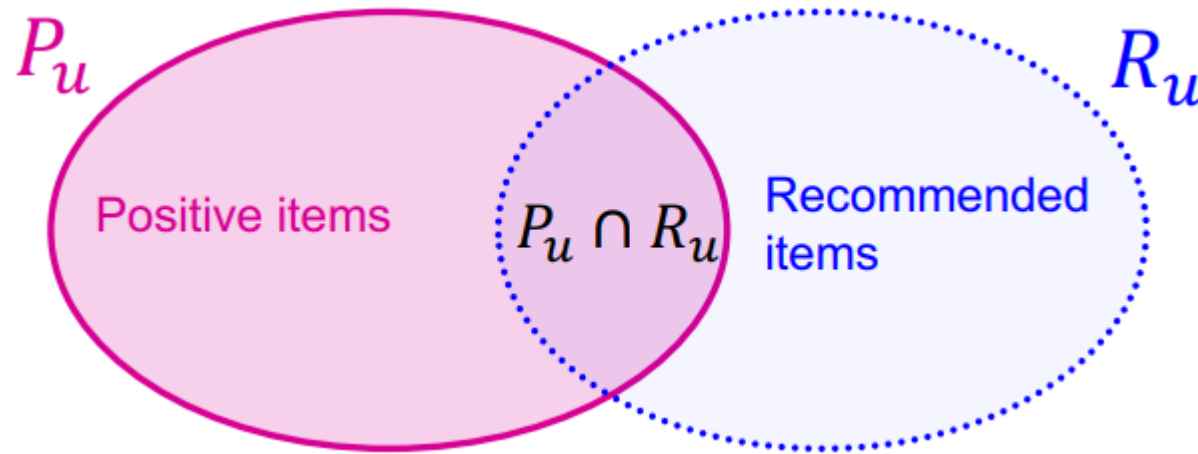


- For each user, we recommend K items.
 - For recommendation to be effective, K needs to be much smaller than the total number of items (up to billions)
 - K is typically in the order of 10—100.
- The goal is to include as many positive items as possible in the top- K recommended items.
 - Positive items = Items that the user will interact with in the future.
- Evaluation metric: Recall@ K

- For each user u :
 - Let P_u be a set of positive items the user will interact in the future.
- Let R_u be a set of items recommended by the model.
 - In top-K recommendation, $|R_u| = K$.
 - Items that the user has already interacted are excluded.



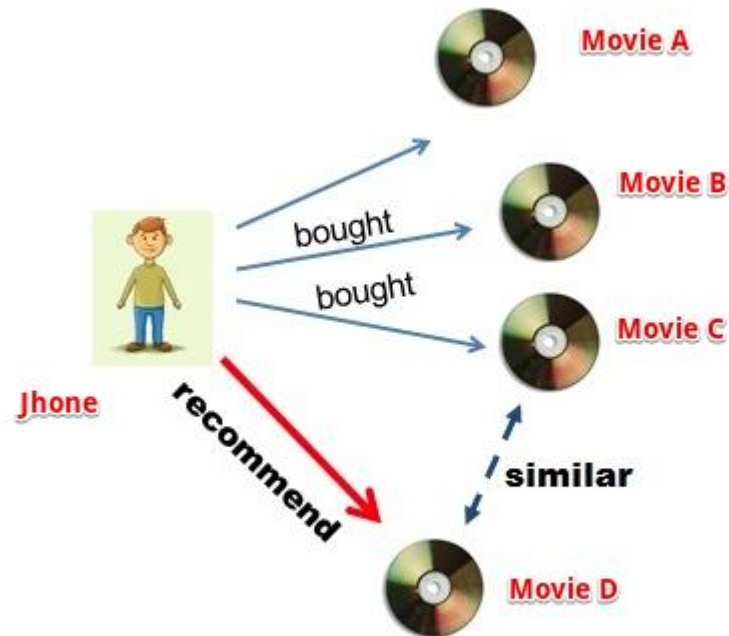
- Recall@K for user u is $|\mathbf{P}_u \cap \mathbf{R}_u|/|\mathbf{P}_u|$
- Higher value indicates more positive items are recommended in top-K for user u .



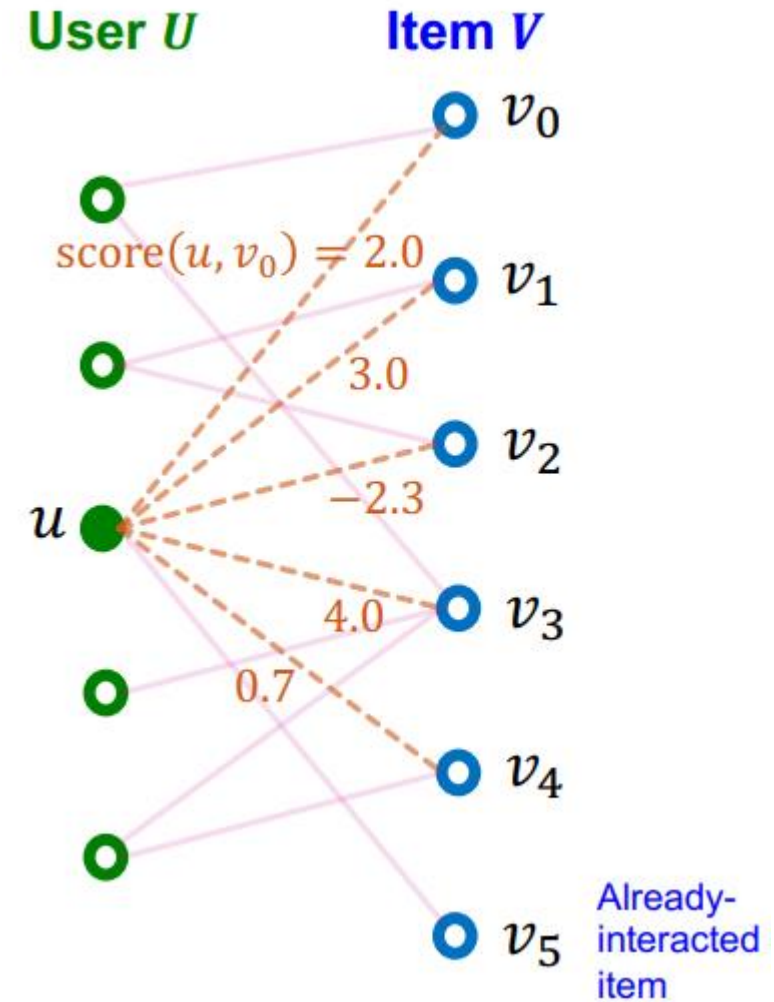
- The final Recall@ K is computed by averaging the recall values across all users.

- U : A set of all users
- V : A set of all items
- E : A set of observed user-item interactions

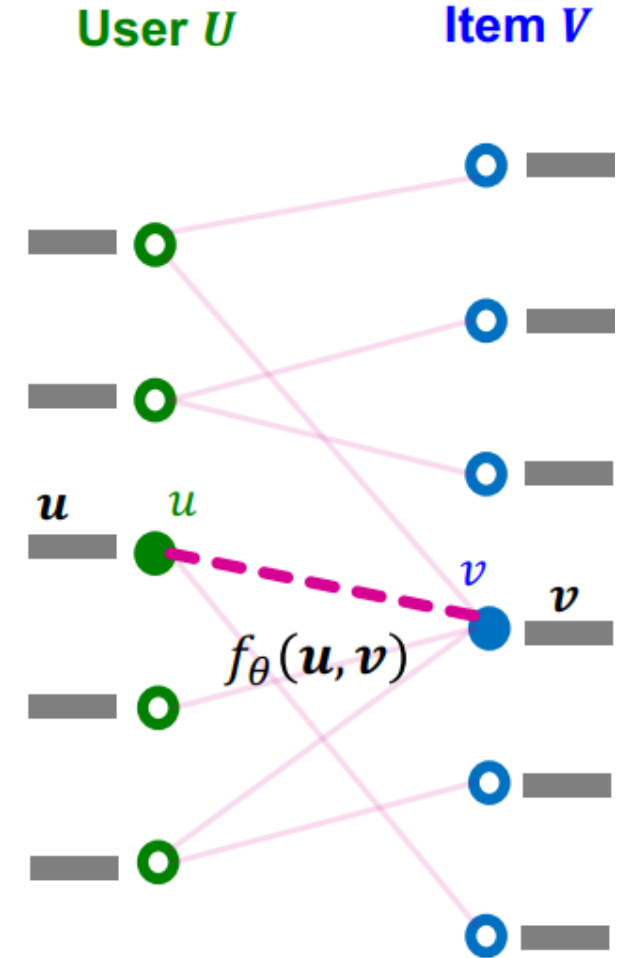
$$E = \{(u, v) \mid u \in U, v \in V, u \text{ interacted with } v\}$$



- To get the top-K items, we need a score function for user-item interaction:
 - For $u \in U, v \in V$, we need to get a real-valued scalar $\text{score}(u, v)$.
 - K items with the largest scores for a given user u (excluding already interacted items) are then recommended
 - i.e., For $K = 2$, recommended items for user u would be v_1, v_3 .



- We consider **embedding-based models** for scoring user-item interactions.
 - For each user $u \in U$, let $u \in R^d$ be its d -dimensional embedding.
 - For each item $v \in V$, let $v \in R^d$ be its d -dimensional embedding.
 - Let $f_\theta(\cdot, \cdot): R^d \times R^d \rightarrow R$ be a parametrized function.
 - Then, $\text{score}(u, v) \equiv f_\theta(u, v)$



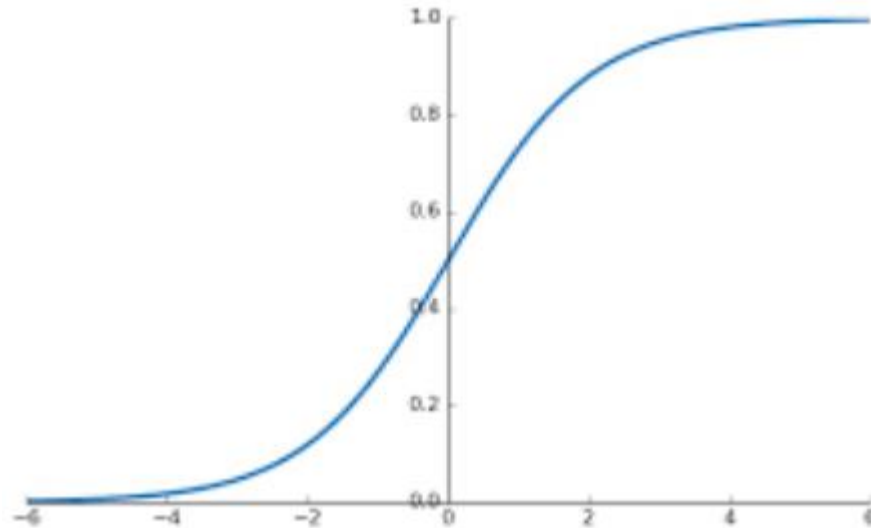
- Embedding-based models have three kinds of parameters:
 - An encoder to generate user embeddings $\{u\}$
 - An encoder to generate item embeddings $\{v\}$
 - Score function $f_{\theta}(\cdot, \cdot)$
- Training objective: Optimize the model parameters to achieve high recall@ K on seen (i.e., training) user-item interactions
 - We hope this objective would lead to high recall@ K on unseen (i.e., test) interactions

- The original training objective ($\text{recall}@K$) is not differentiable.
 - Cannot apply efficient gradient-based optimization.
- Two surrogate loss functions are widely-used to enable efficient gradient-based optimization.
 - Binary loss
 - Bayesian Personalized Ranking (BPR) loss
- Surrogate losses are differentiable and should align well with the original training objective.

- Define positive/negative edges
- A set of **positive edges** E (i.e., observed/training user-item interactions)
- A set of **negative edges** $E_{neg} = \{(u, v) | (u, v) \notin E, u \in U, v \in V\}$
- Define sigmoid function

$$\sigma(x) \equiv \frac{1}{1 + \exp(-x)}$$

- Maps real-valued scores into binary likelihood scores, i.e., in the range of $[0, 1]$.



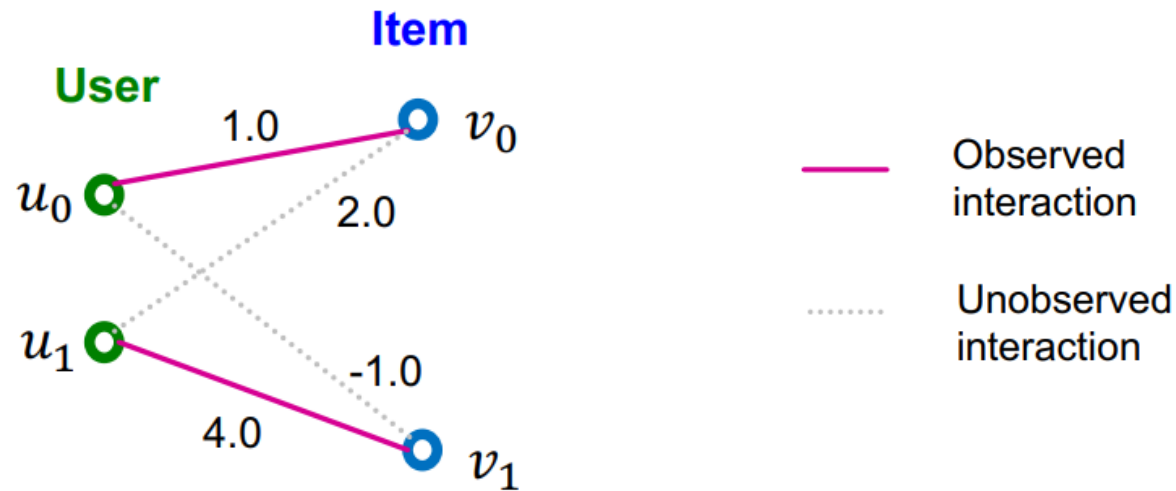
- Binary loss: Binary classification of positive/negative edges using $\sigma(f_\theta(\mathbf{u}, \mathbf{v}))$:

$$-\underbrace{\frac{1}{|E|} \sum_{(u,v) \in E} \log(\sigma(f_\theta(\mathbf{u}, \mathbf{v})))}_{\text{positive edges}} - \underbrace{\frac{1}{|E_{\text{neg}}|} \sum_{(u,v) \in E_{\text{neg}}} \log(1 - \sigma(f_\theta(\mathbf{u}, \mathbf{v})))}_{\text{negative edges}}$$

During training, these terms can be approximated using mini-batch of positive/negative edges

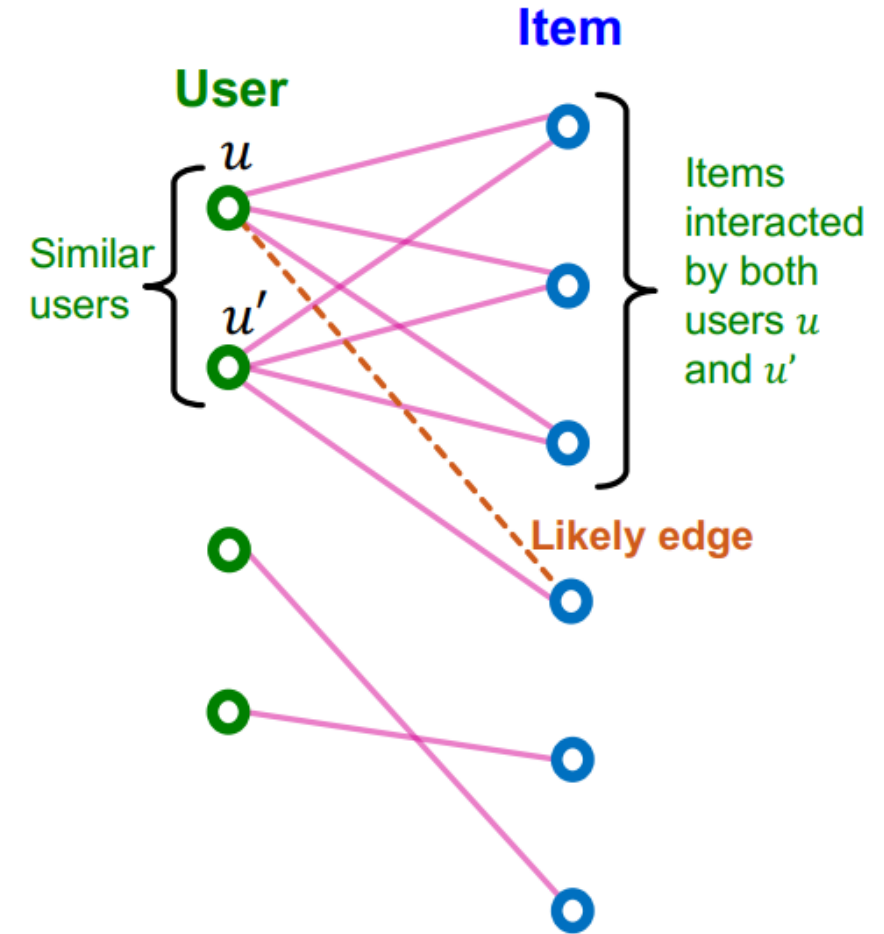
- Binary loss pushes the scores of positive edges higher than those of negative edges.
 - This aligns with the training recall metric since positive edges need to be recalled.

- Surrogate loss function should be defined in a personalized manner.
 - For each user, we want the scores of positive items to be higher than those of the negative items
 - We do not care about the score ordering across users.
- Bayesian Personalized Ranking (BPR) loss achieves this!

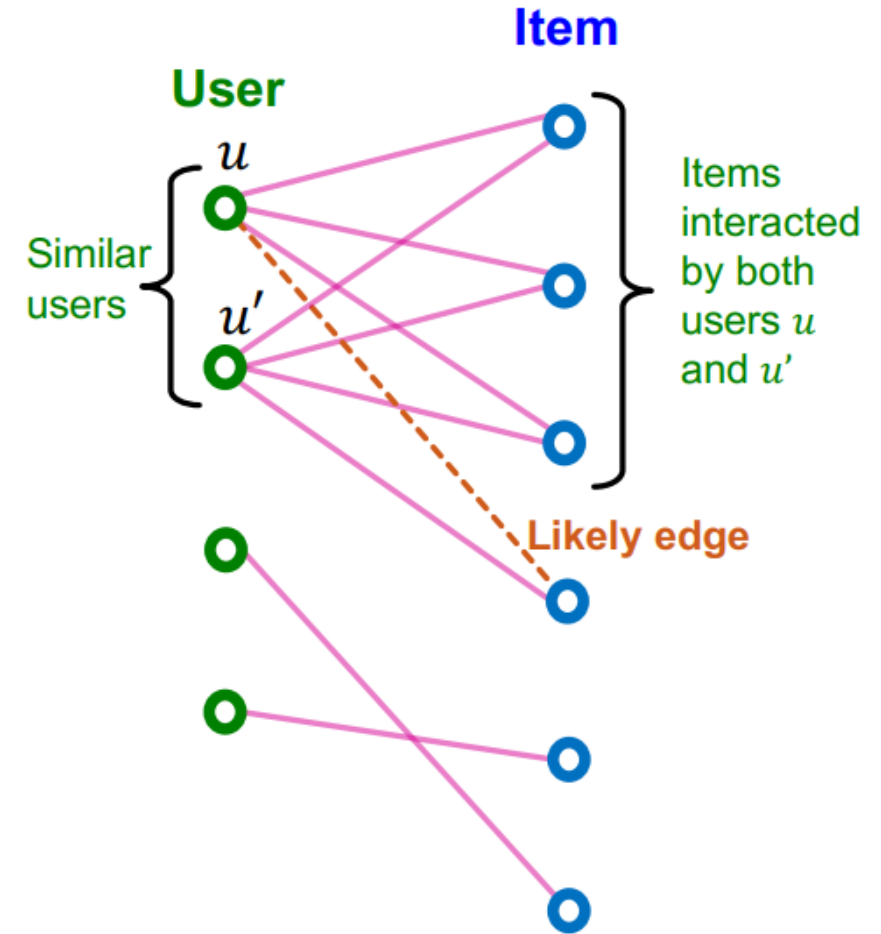


- We have introduced
 - Recall@ K as a metric for personalized recommendation
 - Embedding-based models
 - Three kinds of parameters to learn
 - user encoder to generate user embeddings
 - item encoder to generate item embeddings
 - score function to predict the user-item interaction likelihood.
 - Surrogate loss functions to achieve the high recall metric.
- Embedding-based models have achieved SoTA in recommender systems.
 - Why do they work so well?

- **Underlying idea:** Collaborative filtering
 - Recommend items for a user by collecting preferences of many other similar users.
 - Similar users tend to prefer similar items.
- **Key question:** How to capture similarity between users/items?



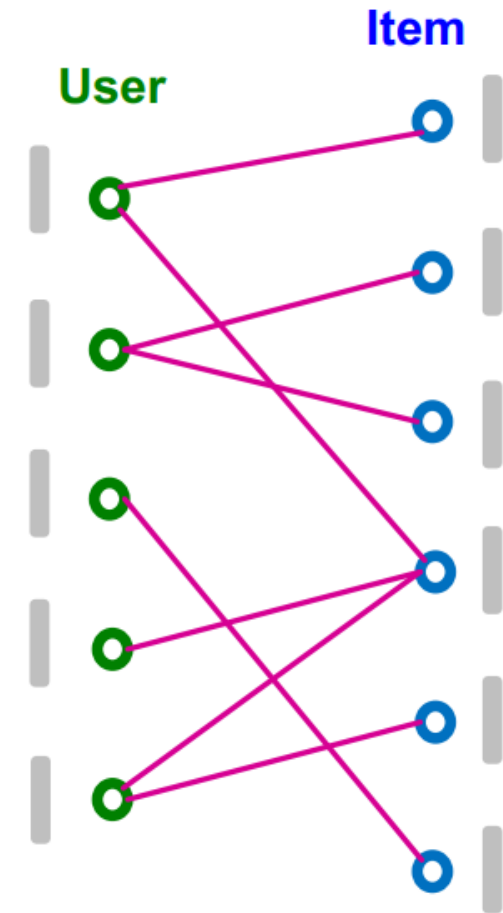
- Embedding-based models can capture similarity of users/items
 - Low-dimensional embeddings cannot simply memorize all user-item interaction data.
 - Embeddings are forced to capture similarity between users/items to fit the data.
 - This allows the models to make effective prediction on unseen user-item interactions.



- Two representative GNN approaches for recommender systems:
 - (1) Neural Graph Collab. Filtering (NGCF) [Wang et al. 2019]
 - (2) LightGCN [He et al. 2020]
 - Improve the conventional collaborative filtering models (i.e., shallow encoders) by explicitly modeling graph structure using GNNs.
 - Assumes no user/item features.
 - (3) KGAT [Wang et al. 2019]

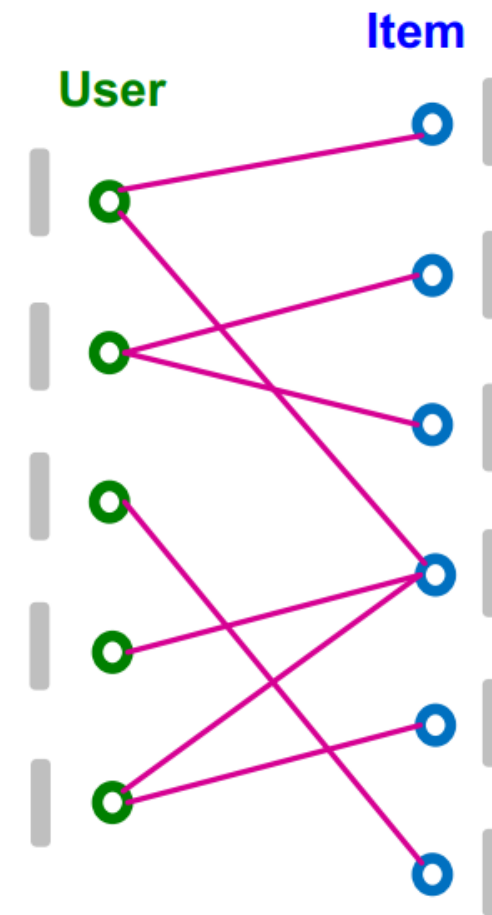
- Conventional collaborative filtering model is based on shallow encoders:
 - No user/item features.
 - Use shallow encoders for users and items
- Score function for user u and item i is:

$$f_{\theta}(\mathbf{u}, \mathbf{v}) \equiv \mathbf{z}_u^T \mathbf{z}_v.$$



Learnable shallow
user/item embeddings

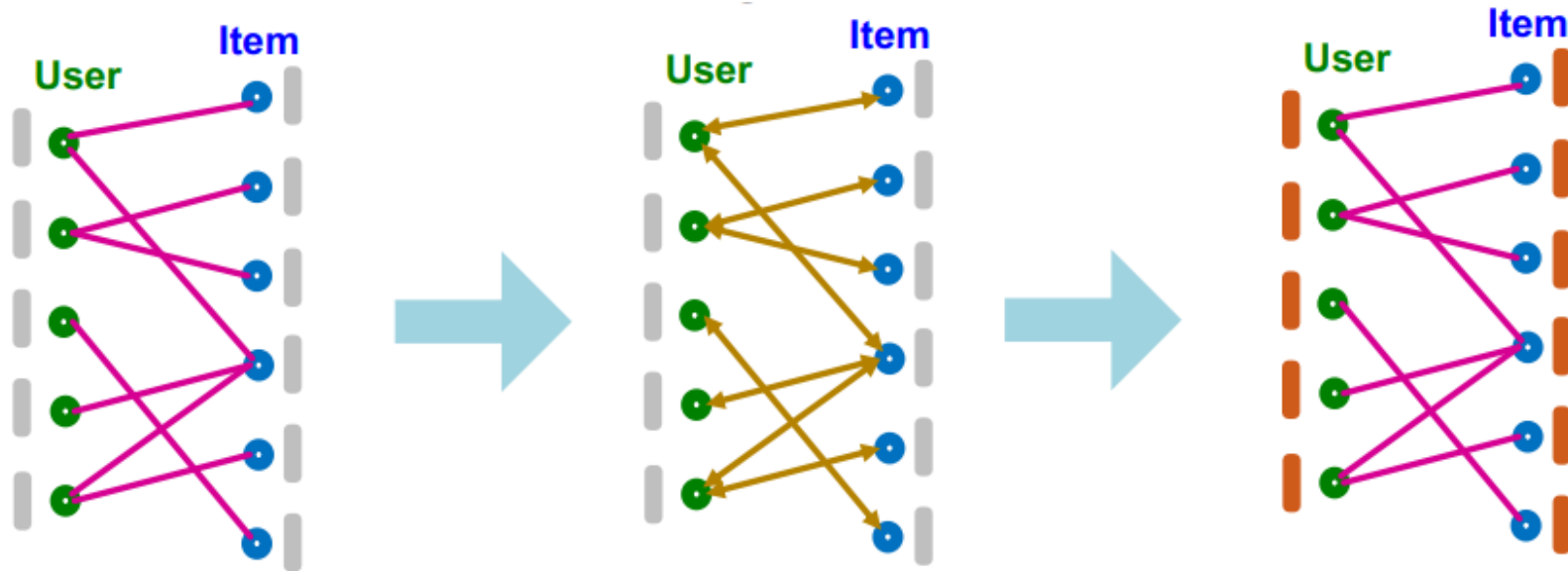
- The model itself does not explicitly capture graph structure
 - The graph structure is only implicitly captured in the training objective.
 - Only the first-order graph structure (i.e., edges) is captured in the training objective.
 - High-order graph structure (e.g., K -hop paths between two nodes) is not explicitly captured.



Learnable shallow
user/item embeddings

- We want a model that:
 - explicitly captures graph structure (beyond implicitly through the training objective)
 - captures high-order graph structure (beyond the first-order edge connectivity structure)
- GNNs are a natural approach to achieve both
 - Neural Graph Collaborative Filtering (NGCF) [Wang et al. 2019]
 - LightGCN [He et al. 2020]
 - A simplified and improved version of NGCF

- Neural Graph Collaborative Filtering (NGCF) explicitly incorporates high-order graph structure when generating user/item embeddings.
- **Key idea:** Use a GNN to generate graph-aware user/item embeddings.

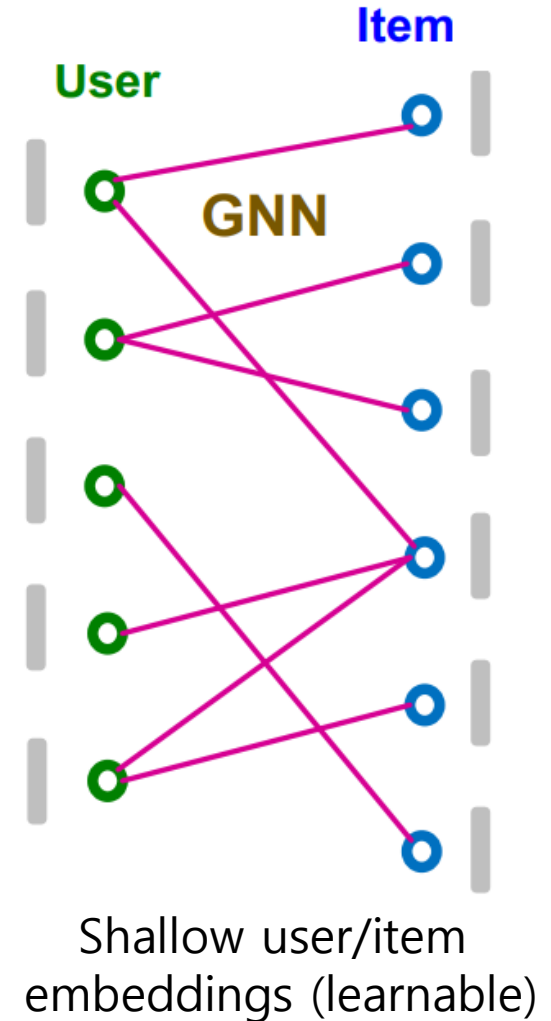


Initial shallow embeddings
(not graph-aware)

Use a GNN to
propagate embeddings

NGCF's graph-aware
embeddings

- Given: User-item bipartite graph.
- NGCF framework:
 - Prepare shallow learnable embedding for each node.
 - Use multi-layer GNNs to propagate embeddings along the bipartite graph.
 - High-order graph structure is captured.
 - Final embeddings are explicitly graphaware
- Two kinds of learnable params are jointly learned:
 - Shallow user/item embeddings
 - GNN's parameters

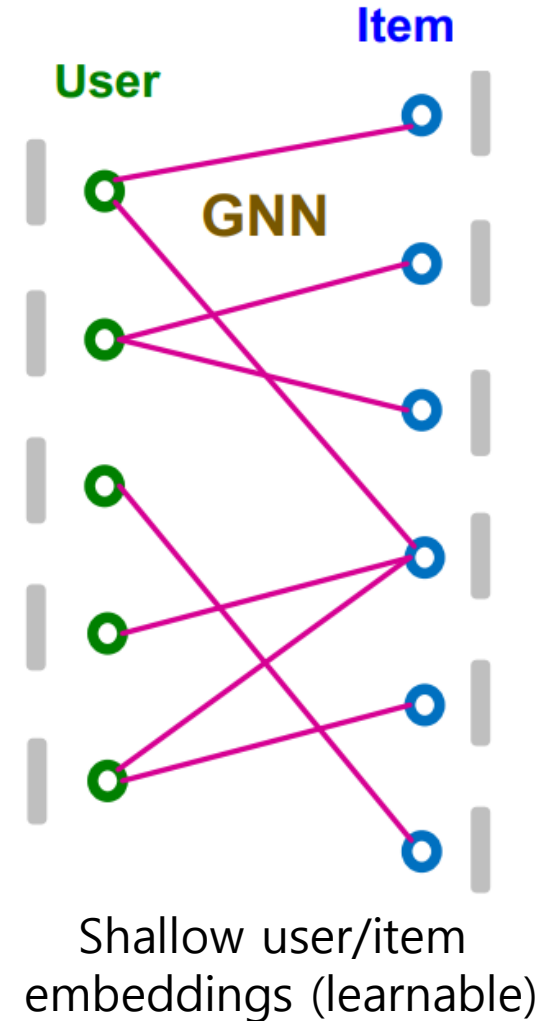


- Set the shallow learnable embeddings as the initial node features:
 - For every user $u \in U$, set h_u^0 as the user's shallow embedding.
 - For every item $v \in V$, set h_v^0 as the item's shallow embedding.
- **Neighbour aggregation:**
 - Iteratively update node embeddings using neighboring embeddings.

$$\mathbf{h}_v^{(k+1)} = \text{COMBINE} \left(\mathbf{h}_v^{(k)}, \text{AGGR} \left(\left\{ \mathbf{h}_u^{(k)} \right\}_{u \in N(v)} \right) \right)$$

$$\mathbf{h}_u^{(k+1)} = \text{COMBINE} \left(\mathbf{h}_u^{(k)}, \text{AGGR} \left(\left\{ \mathbf{h}_v^{(k)} \right\}_{v \in N(u)} \right) \right)$$

- High-order graph structure is captured through iterative neighbor aggregation.

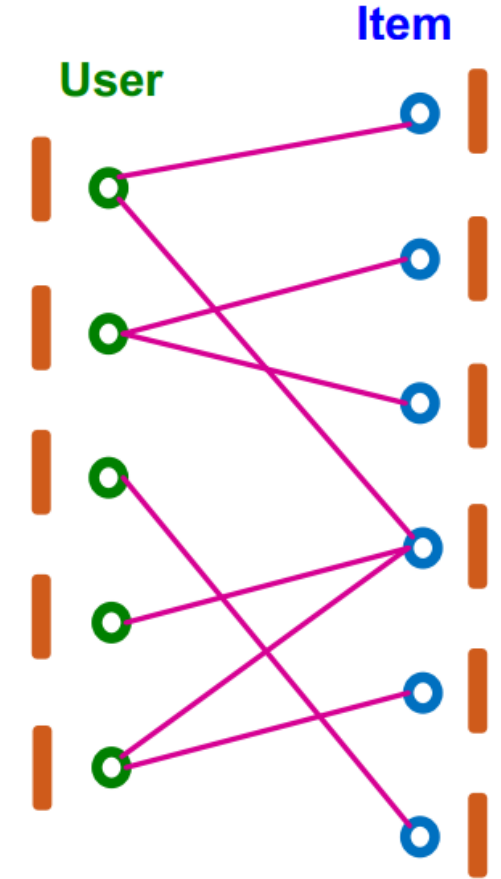


- After k rounds of neighbor aggregation, we get the final user/item embeddings $h_u^{(k)}$ and $h_v^{(k)}$.
- we set

$$\mathbf{u} \leftarrow \mathbf{h}_u^{(K)}, \mathbf{v} \leftarrow \mathbf{h}_v^{(K)}.$$

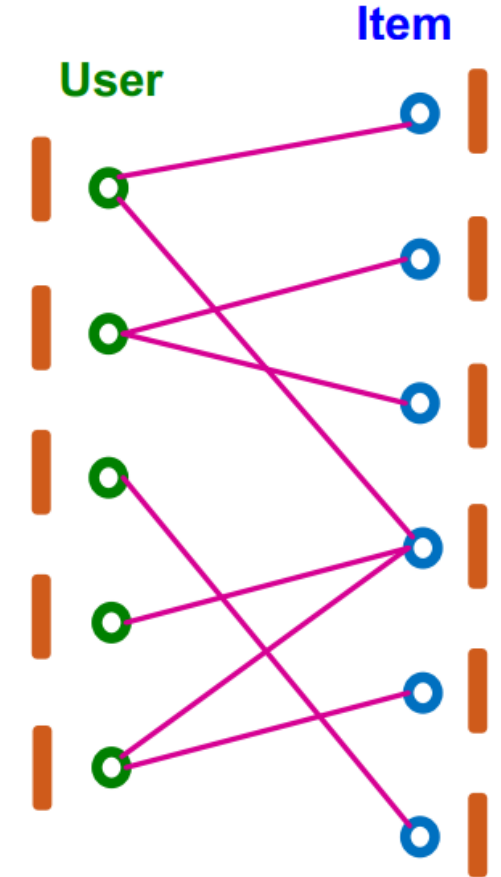
- Score function is the inner product

$$\text{score}(u, v) = \mathbf{u}^T \mathbf{v}$$



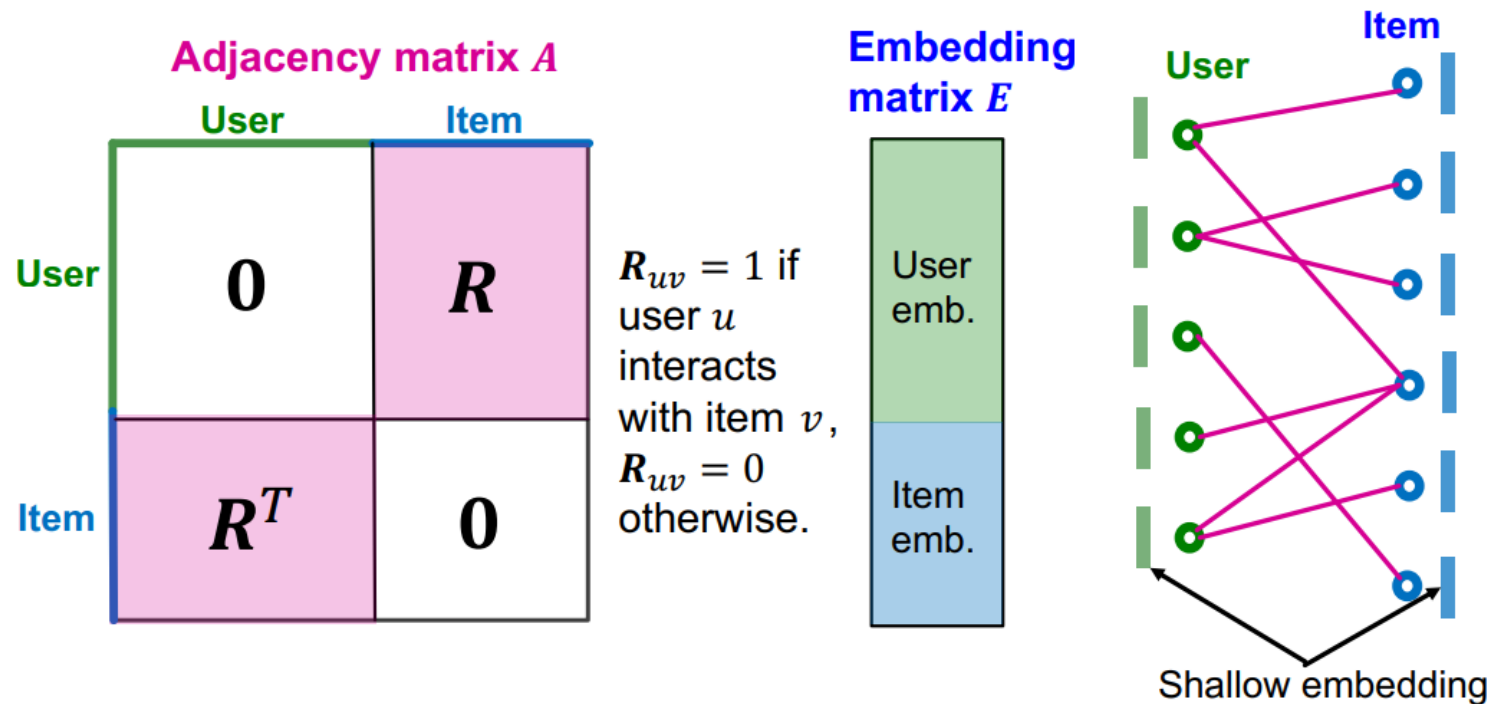
Final user/item embeddings (graph-aware)

- Conventional collaborative filtering uses shallow user/item embeddings.
 - The embeddings do not explicitly model graph structure.
 - The training objective does not model high-order graph structure.
- NGCF uses a GNN to propagate the shallow embeddings.
 - The embeddings are explicitly aware of high-order graph structure.

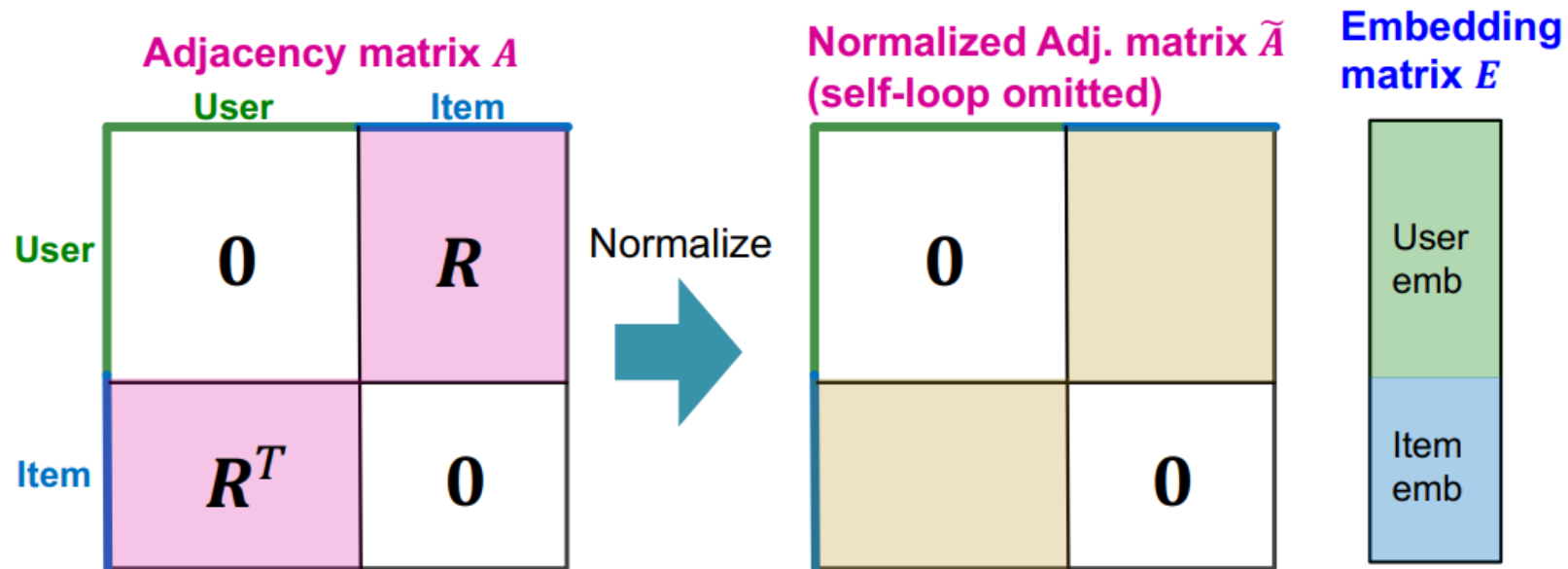


- **Recall:** NGCF jointly learns two kinds of parameters:
 - Shallow user/item embeddings (quite expressive)
 - GNN's parameters
- **Can we simplify the GNN used in NGCF (e.g., remove its learnable parameters)?**
- **Overview of the idea:**
 - Adjacency matrix for a bipartite graph
 - Matrix formulation of GCN
 - Simplification of GCN by removing non-linearity

- Adjacency matrix of a (undirected) bipartite graph.
- Shallow embedding matrix.



- Given:
 - Adjacency matrix A
 - Initial learnable embedding matrix E



- **Define:** The diffusion matrix
- Let D be the degree matrix of A .
- Define the normalized adjacency matrix \tilde{A} as

$$\tilde{A} \equiv D^{-1/2} A D^{-1/2}$$

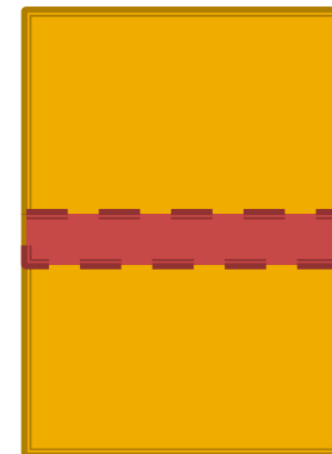
- Let E^k be the embedding matrix at k -th layer.
- Each layer of GCN's aggregation can be written in a matrix form:

$$E^{(k+1)} = \text{ReLU}(\tilde{A}E^{(k)}W^{(k)})$$

Neighbor aggregation

Learnable linear transformation

Matrix of node embeddings $E^{(k)}$



Each row stores node embedding

- Removing ReLU significantly simplifies GCN

$$\mathbf{E}^{(K)} = \boxed{\tilde{\mathbf{A}}^K \mathbf{E}} \mathbf{W} \quad \mathbf{W} \equiv \mathbf{W}^{(0)} \dots \mathbf{W}^{(K-1)}$$

- Diffusing node embeddings along the graph
- Algorithm: Apply $\mathbf{E} \leftarrow \tilde{\mathbf{A}} \mathbf{E}$ for K times
- Each matrix multiplication diffuses the current embeddings to their one-hop neighbors.
- Note: $\tilde{\mathbf{A}}^K$ is dense and never gets materialized. Instead, the above iterative matrix-vector product is used to compute $\tilde{\mathbf{A}}^K \mathbf{E}$.

- The embedding propagation of LightGCN is closely related to GCN
- Recall: GCN (neighbor aggregation part)

$$\mathbf{h}_v^{(k+1)} = \sum_{u \in N(v)} \frac{1}{\sqrt{d_u} \sqrt{d_v}} \cdot \mathbf{h}_u^{(k)}$$

Node degree

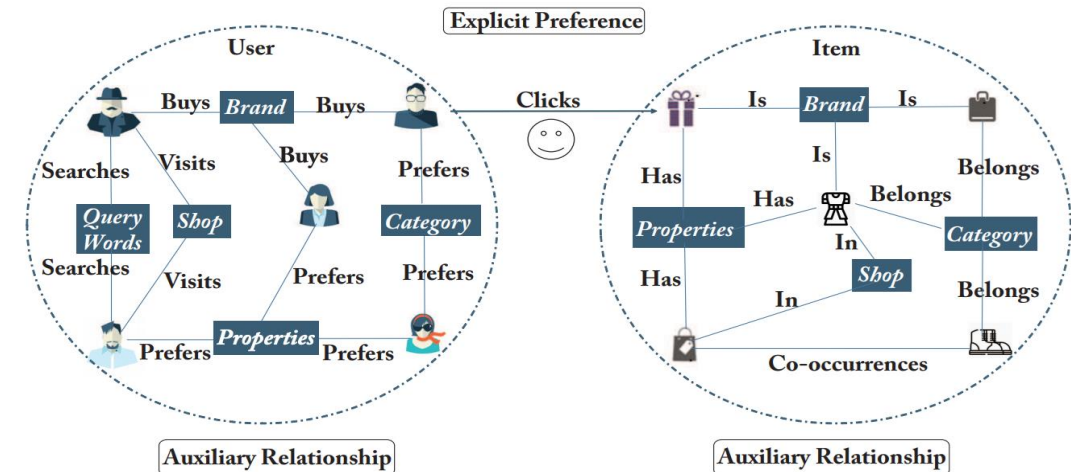
- Self-loop is added in the neighborhood definition.
- LightGCN uses the same equation except that
 - Self-loop is not added in the neighborhood definition.
 - Final embedding takes the average of embeddings from all the layers:

$$\mathbf{h}_v = \frac{1}{K+1} \sum_{k=0}^K \mathbf{h}_v^{(k)}.$$

- LightGCN simplifies NGCF by removing the learnable parameters of GNNs.
- Learnable parameters are all in the shallow input node embeddings.
 - Diffusion propagation only involves matrix-vector multiplication.
 - The simplification leads to better empirical performance than NGCF.

➤ Knowledge graph:

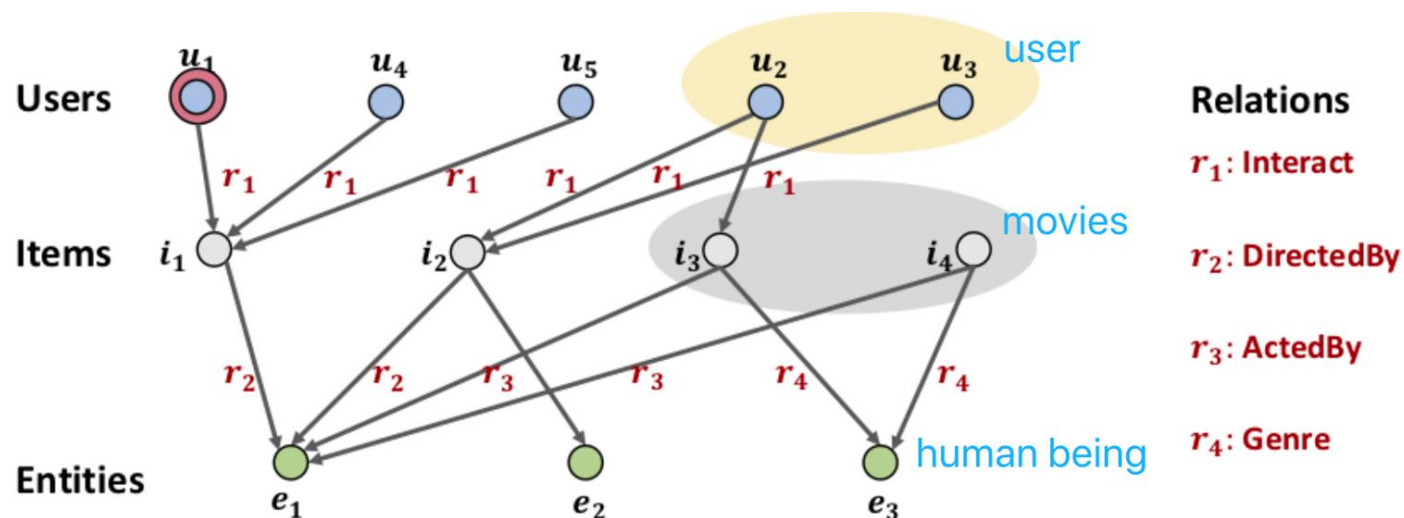
- Rich semantic relatedness among items in KG
- Improve the accuracy of prediction result
- Can reasonable extend a user's interest



➤ Difficulty: How to use KG to enhance the prediction?

- Integration: Treat two graphs separately or merge them?
- Graph Simplification: Do we make all efforts on the KG, or on its sub-graph?
- Multi-relation Propagation: There are different relationship. Should we give them the same weight for a certain user?

- KGAT:
 - Combine two graphs into one
 - Use attention mechanism by (entity1, relation, entity2)

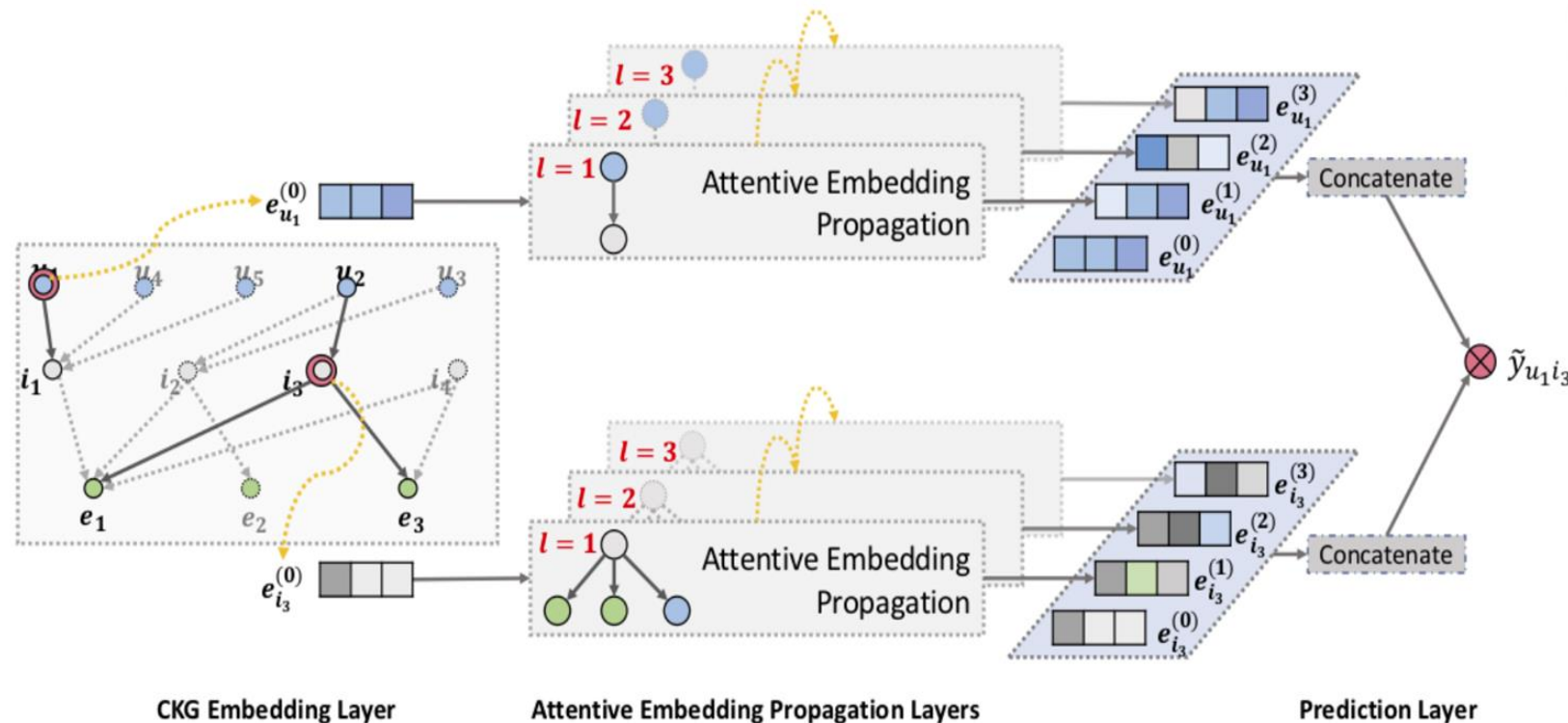


➤ KGAT:

- For each item and user node, use attention to gather their neighbor's information
- Concatenate embeddings from different jump ($l=1, l=2, l=3 \dots l=L$)

$$\pi_r^u = g(r, u)$$

$$v_{N(v)}^u = \sum_{e \in N(v)} \pi_{r_v, e}^u$$





네트워크 과학연구실
NETWORK SCIENCE LAB



가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

